

## HW-5

Q1a)

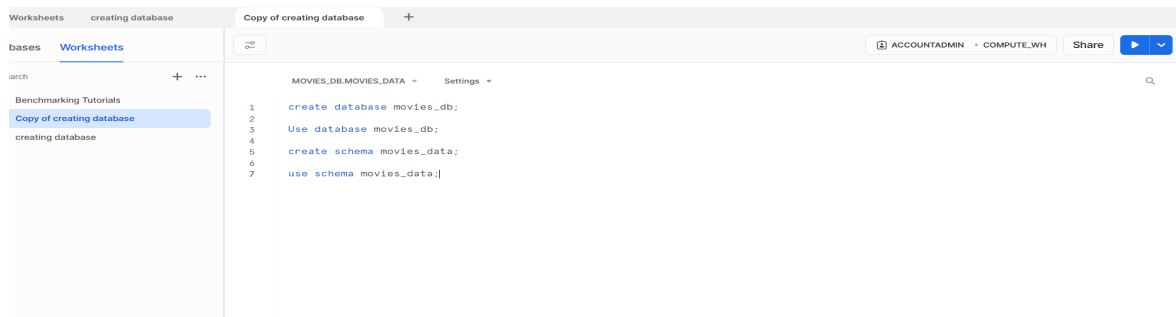
I am using Snowflake in the first part. I did not use it in my project that is why I explored Snowflake for this question.

The dataset I used is a movie dataset. Below is the drive link for the database.

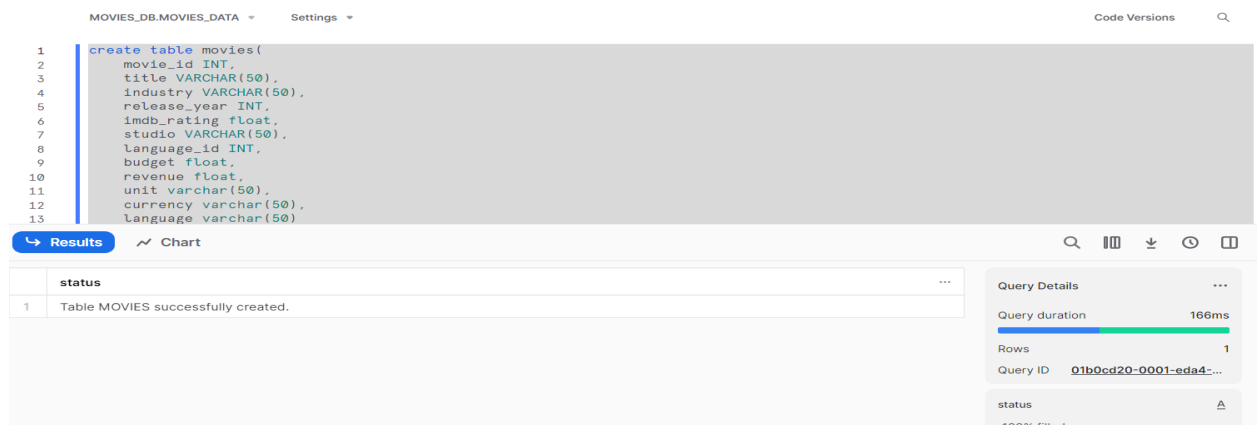
<https://docs.google.com/spreadsheets/d/1JuNNgR50-815QoptzElm4RYZrJE9eQv0/edit?usp=sharing&ouid=105974981587456835079&rtpof=true&sd=true>

- Database creation and schema creation

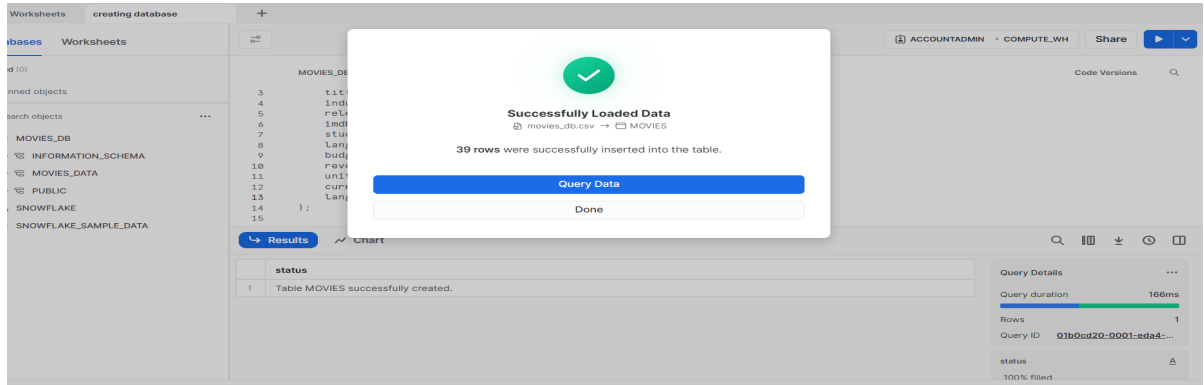
The very first step I performed is database creation. I created a dataset named `movies_db`.



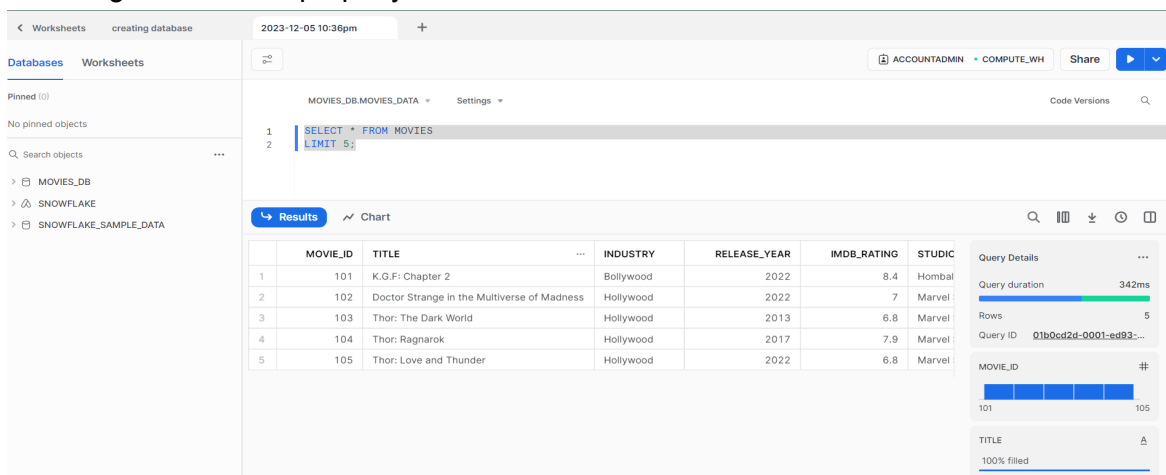
- Created a table to load the dataset which I downloaded



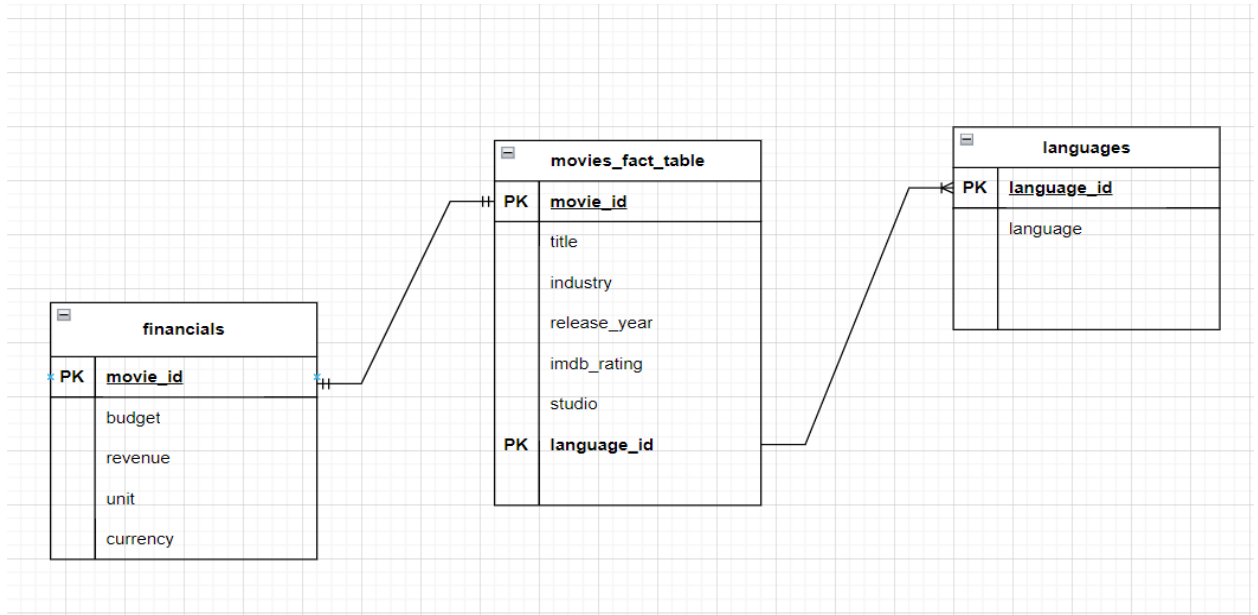
- Loaded the dataset into the movies table which I created and the snippet is above about the creation of movies table.



- Checking if the data is properly loaded or not.



- Below is the star schema of the dataset. It is only possible to create two dimension tables and not more than that. If I create different dimension tables for currency or unit it will only have redundant data. Instead I kept all that information in the financials dimension table.
- Also In languages dimension table I tried to save a lot of space as I mapped a single language\_id to each language and the maximum number of rows required here are only 8, whereas in the main table this data was filled in every row.



- Created dimension table financials according to the star schema which I have created.

MOVIES\_DB.MOVIES\_DATA

```

1 CREATE TABLE financials (
2   MOVIE_ID INT PRIMARY KEY,
3   BUDGET FLOAT,
4   REVENUE FLOAT,
5   UNIT VARCHAR(50),
6   CURRENCY VARCHAR(50)
7 )
  
```

Results

status
1 Table FINANCIALS successfully created.

Query Details

Query duration 475ms

Rows 1

Query ID 01b0cd35-0001-edf5-0...

status 100% filled

- Created another dimension table of languages according to the star schema which I have created.

MOVIES\_DB.MOVIES\_DATA

```

1 CREATE TABLE Languages (
2   language_id INT PRIMARY KEY,
3   language VARCHAR(50)
4 )
  
```

Results

status
1 Table LANGUAGES successfully created.

Query Details

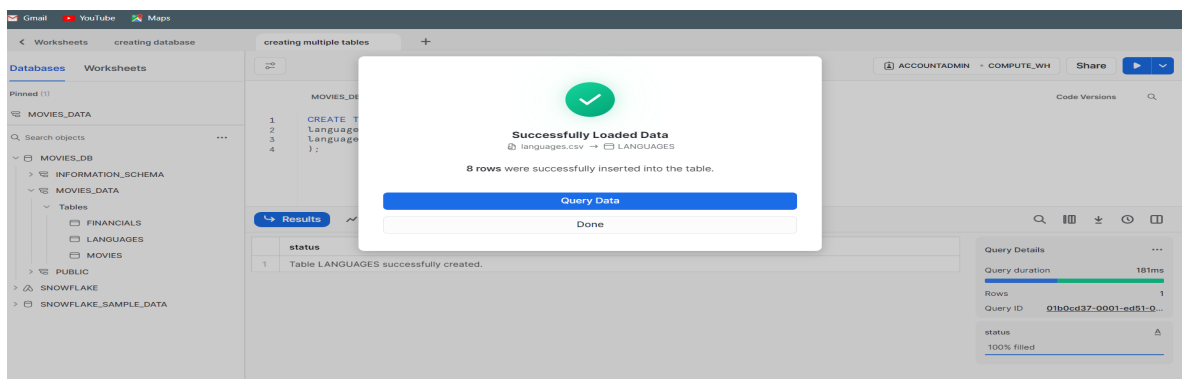
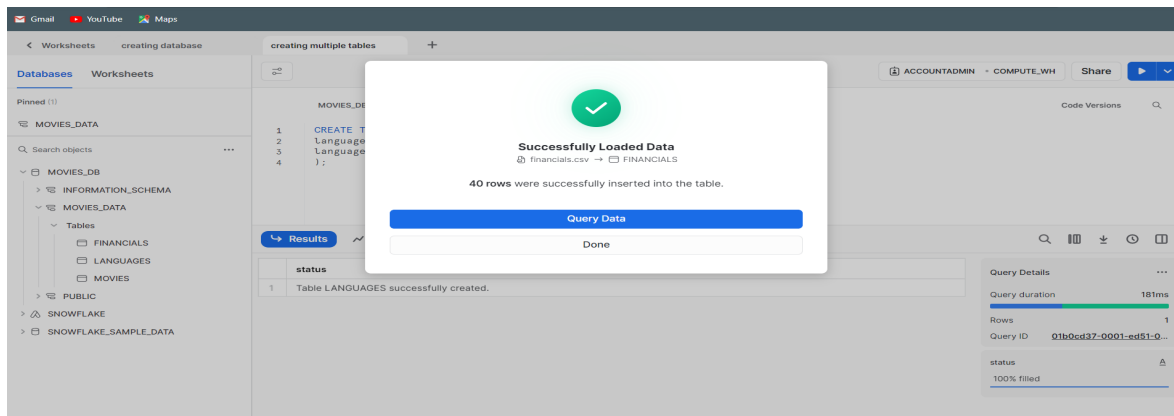
Query duration 181ms

Rows 1

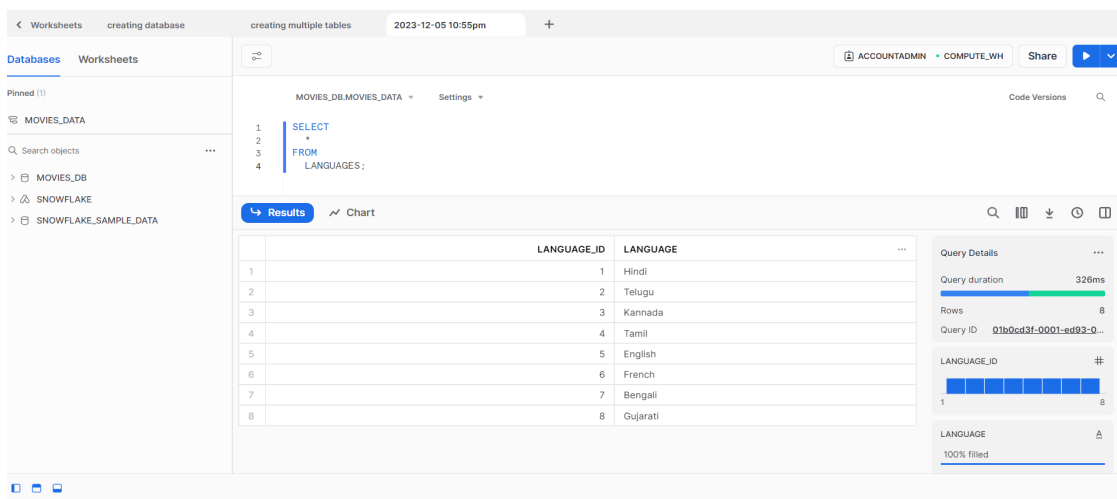
Query ID 01b0cd37-0001-ed51-0...

status 100% filled

- These two below are the screenshots of loading data into newly created dimension tables.



- Below screenshot is just of verification that whether the data has properly been loaded or not.



- Now I have created the fact table where movie\_id is the primary key and it has 2 dimension tables named financials and languages.

The screenshot shows the Snowflake SQL Editor with the following SQL query:

```

1 create table movies_fact_table(
2     movie_id INT,
3     title VARCHAR(50),
4     industry VARCHAR(50),
5     release_year INT,
6     imdb_rating float,
7     studio VARCHAR(50),
8     language_id INT,
9     primary key (movie_id),
10    foreign key (movie_id) references financials(movie_id),
11    foreign key (language_id) references languages(language_id)
12 );

```

The Results tab shows a single row with the status: "Table MOVIES\_FACT\_TABLE successfully created."

Query Details: Query duration 255ms, Rows 1, Query ID 01b0cd49-0001-ed45-...

- Below are the screenshots which verifies that data has been successfully loaded into the fact table.

The screenshot shows a modal dialog titled "Successfully Loaded Data" with a green checkmark icon. The message states: "39 rows were successfully inserted into the table." Below the message are two buttons: "Query Data" and "Done".

The background shows the SQL Editor with the same CREATE TABLE query as the previous screenshot. The Results tab shows the same status message: "Table MOVIES\_FACT\_TABLE successfully created."

Query Details: Query duration 255ms, Rows 1, Query ID 01b0cd49-0001-ed45-...

The screenshot shows the SQL Editor with the following SQL query:

```

1 SELECT *
2 FROM
3     MOVIES_FACT_TABLE
4 LIMIT
5     10;

```

The Results tab displays a table with 10 rows of data:

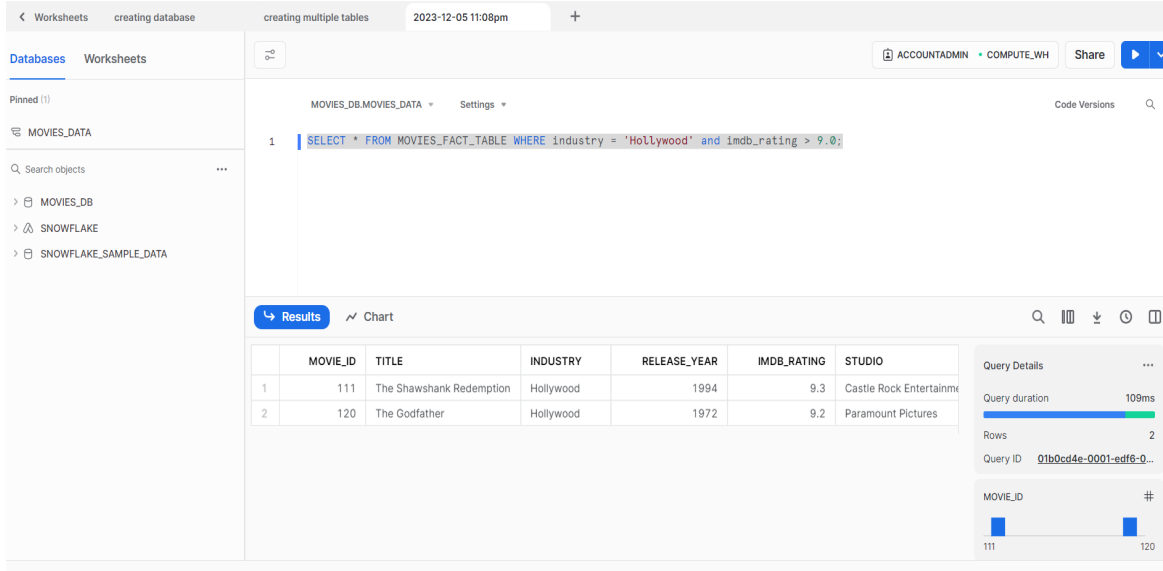
	MOVIE_ID	TITLE	INDUSTRY	RELEASE_YEAR	IMDB_RATING	STUDIO
1	101	K.G.F: Chapter 2	Bollywood	2022	8.4	Homb
2	102	Doctor Strange in the Multiverse of Madness	Hollywood	2022	7	Marve
3	103	Thor: The Dark World	Hollywood	2013	6.8	Marve
4	104	Thor: Ragnarok	Hollywood	2017	7.9	Marve
5	105	Thor: Love and Thunder	Hollywood	2022	6.8	Marve
6	106	Sholay	Bollywood	1975	8.1	Unitec
7	107	Dilwale Dulhania Le Jayenge	Bollywood	1995	8	Yash f
8	108	3 Idiots	Bollywood	2009	8.4	Vinod
9	109	Kabhi Khushi Kabhie Gham	Bollywood	2001	7.4	Dharm
10	110	Bajirao Mastani	Bollywood	2015	7.2	Unitec

Query Details: Query duration 559ms, Rows 10, Query ID 01b0cd4e-0001-ed45-...

A progress bar for MOVIE\_ID shows 100% filled, with a range from 101 to 110.

## Analytics part(Queries)

1) Below is the basic query which I tried to find out about the movies released by hollywood and their imdb rating is more than 9.



The screenshot shows a Snowflake query editor interface. The query is:

```
SELECT * FROM MOVIES_FACT_TABLE WHERE industry = 'Hollywood' and imdb_rating > 9.0;
```

The results table displays the following data:

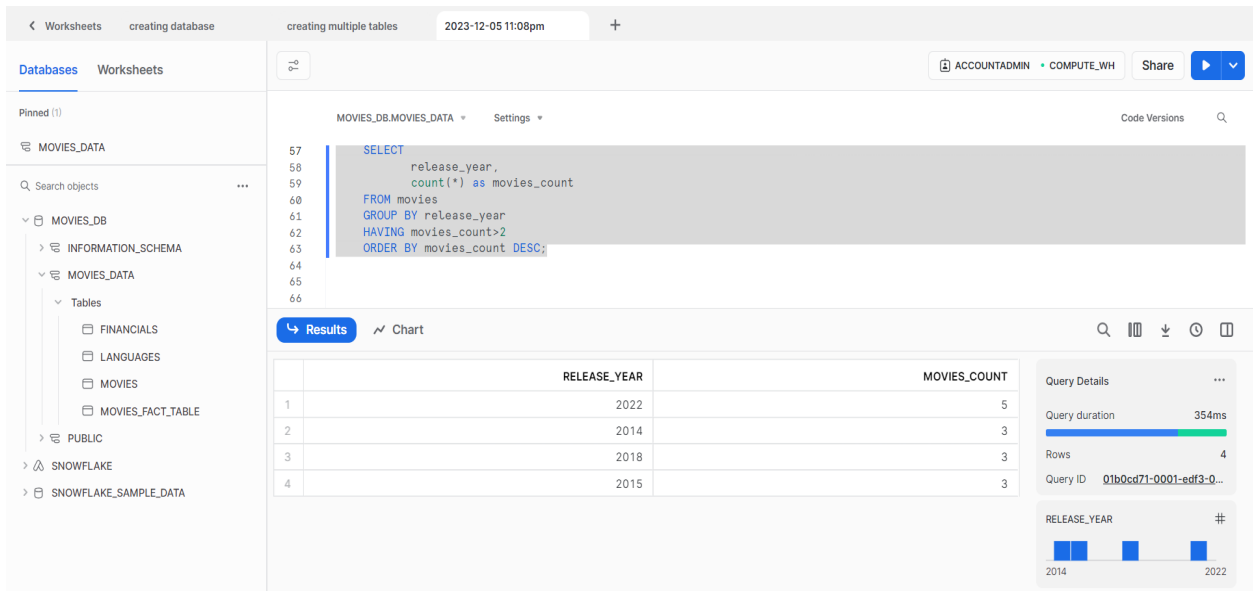
	MOVIE_ID	TITLE	INDUSTRY	RELEASE_YEAR	IMDB_RATING	STUDIO
1	111	The Shawshank Redemption	Hollywood	1994	9.3	Castle Rock Entertainment
2	120	The Godfather	Hollywood	1972	9.2	Paramount Pictures

Query Details:

- Query duration: 109ms
- Rows: 2
- Query ID: 01b0cd4e-0001-edf6-9...

A bar chart shows the distribution of MOVIE\_ID values, with bars for 111 and 120.

2) Below I tried to print all the years where more than 2 movies were released



The screenshot shows a Snowflake query editor interface. The query is:

```
SELECT release_year, count(*) as movies_count FROM movies GROUP BY release_year HAVING movies_count > 2 ORDER BY movies_count DESC;
```

The results table displays the following data:

	RELEASE_YEAR	MOVIES_COUNT
1	2022	5
2	2014	3
3	2018	3
4	2015	3

Query Details:

- Query duration: 354ms
- Rows: 4
- Query ID: 01b0cd71-0001-edf3-0...

A bar chart shows the distribution of RELEASE\_YEAR values, with bars for 2014, 2015, 2018, and 2022.

3)Below I joined all the three tables and printed the movies whose currency is 'INR' and Unit is 'Millions' and the budget is more than 500 where the budget is calculated on the basis of unit unit.

MOVIES\_DB.MOVIES\_DATA

```

1 SELECT m.title,m.imdb_rating, f.budget, f.unit, f.currency, l.language
2 FROM MOVIES_FACT_TABLE m
3 JOIN languages l
4 ON m.language_id = l.language_id
5 JOIN financials f
6 ON m.movie_id = f.movie_id
7 where currency = 'INR' and unit = 'Millions' and budget>500;

```

	TITLE	IMDB_RATING	BUDGET	UNIT	CURRENCY	LANGUAGE
1	3 Idiots	8.4	550	Millions	INR	Hindi
2	PK	8.1	850	Millions	INR	Hindi
3	Bajrangi Bhaijaan	8.1	900	Millions	INR	Hindi

Query Details: Query duration 83ms, Rows 3, Query ID 01b0cd5c-0001-edf5-0...

4)Below I have used Common Table Expression(CTE's) and printed movies that produced 500% profit.

MOVIES\_DB.MOVIES\_DATA

```

8
9
10 with profit_table as (
11     select
12         movie_id as id,
13         ((revenue-budget)/budget)* 100 as profit_pct,
14         unit AS unit
15     from financials
16 )
17 select P.id, P.profit_pct, M.TITLE, P.unit from profit_table p
18 JOIN movies_fact_table m
19 ON p.id = m.movie_id
20 where profit_pct>500;

```

	ID	PROFIT_PCT	TITLE	UNIT
1	101	1150	K.G.F: Chapter 2	Billions
2	108	627.272727273	3 Idiots	Millions
3	117	1001	Titanic	Millions
4	119	1101.265822785	Avatar	Millions
5	120	3941.666666667	The Godfather	Millions
6	122	1364.545454545	Schindler's List	Millions

Query Details: Query duration 217ms, Rows 13, Query ID 01b0cd62-0001-edf3-0...

5)Below I have printed movies that produced 500% profit and their rating was less than average rating for all movies.

MOVIES\_DB.MOVIES\_DATA

```

19
20
21
22
23 select
24     x.movie_id, y.title, x.pct_profit, y.imdb_rating
25 from (
26     select
27         *,
28         (revenue-budget)*100/budget as pct_profit
29     from financials
30 ) x
31 join (select * from movies where imdb_rating < (select avg(imdb_rating) from movies)) y
32 on x.movie_id=y.movie_id
33 where pct_profit>500;
34

```

	MOVIE_ID	TITLE	PCT_PROFIT	IMDB_RATING
1	119	Avatar	1101.265822785	7.8

Query Details: Query duration 365ms, Rows 1, Query ID 01b0cd6a-0001-edf6-0...

6) Below I have printed Movies with highest and lowest ratings from the dataset.

The screenshot shows the Snowflake web interface. On the left, the 'Databases' sidebar is open, showing the hierarchy: MOVIES\_DB > MOVIES\_DATA > TABLES > MOVIES. The main editor displays the following SQL query:

```
select * from movies where imdb_rating in (
(select min(imdb_rating) from movies),
(select max(imdb_rating) from movies)
);
```

The 'Results' tab shows a table with 2 rows:

	MOVIE_ID	TITLE	INDUSTRY	RELEASE_YEAR	IMDB_RATING	STUDIO
1	111	The Shawshank Redemption	Hollywood	1994	9.3	Castle Rock Entertainment
2	139	Race 3	Bollywood	2018	1.9	Salman Khan Films

Query Details on the right: Query duration 337ms, Rows 2, Query ID 01b0cd8c-0001-edf6-9...

7) Below I have printed all the movies with minimum and maximum release\_year.

The screenshot shows the Snowflake web interface. The main editor displays the following SQL query:

```
select * from movies where release_year in (
(select min(release_year) from movies),
(select max(release_year) from movies)
)
```

The 'Results' tab shows a table with 6 rows:

	MOVIE_ID	TITLE	INDUSTRY	RELEASE_YEAR	IMDB_RATING	STUDIO
1	101	K.G.F: Chapter 2	Bollywood	2022	8.4	Hombale Films
2	102	Doctor Strange in the Multiverse of Madness	Hollywood	2022	7	Marvel Studios
3	105	Thor: Love and Thunder	Hollywood	2022	6.8	Marvel Studios
4	118	It's a Wonderful Life	Hollywood	1946	8.6	Liberty Films
5	133	RRR	Bollywood	2022	8	DVVFilms
6	135	The Kashmir Files	Bollywood	2022	8.3	Zee Studios

Query Details on the right: Query duration 81ms, Rows 6, Query ID 01b0cd8e-0001-edf6-9...

8) Below I have printed the average rating of movies per studio and also ordered them by average rating in descending format.

The screenshot shows the Snowflake web interface. The main editor displays the following SQL query:

```
SELECT
studio,
count(studio) as cnt,
round(avg(imdb_rating), 1) as avg_rating
from movies WHERE studio != ''
GROUP BY studio
order by avg_rating DESC;
```

The 'Results' tab shows a table with 8 rows:

	STUDIO	CNT	AVG_RATING
1	Castle Rock Entertainment	1	9.3
2	Syncopy	1	9
3	Warner Bros. Pictures	2	8.7
4	Paramount Pictures	2	8.6
5	Universal Pictures	2	8.6
6	Liberty Films	1	8.6
7	Universal Pictures	1	8.5
8	Hombale Films	1	8.4

Query Details on the right: Query duration 386ms, Rows 22, Query ID 01b0cd70-0001-edd4-...

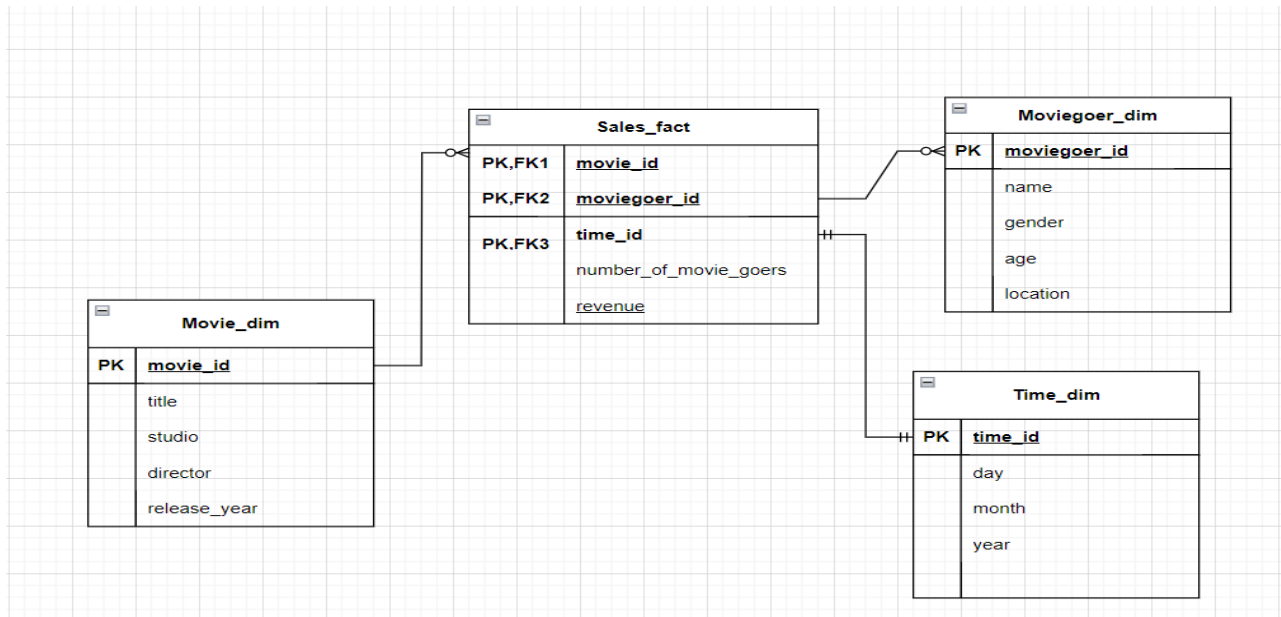


Q1)b) Below is the link to the blog written on Medium:

<https://medium.com/@yeolevishal11/insights-from-a-movies-dataset-using-snowflake-d50b809b499d>

Q2)

a)



b) Total revenue collected through ticket sales by each movie in the pandemic year, 2020.

```
SELECT movie_dim.title, SUM(Sales_fact.revenue) AS total_revenue
FROM Sales_fact s
JOIN movie_dim m
ON s.movie_id = m.movie_id
JOIN Time_dim t
ON s.time_id = t.time_id
WHERE t.year = 2020
GROUP BY m.title;
```

c)

```
SELECT movie, SUM(ticket_price) AS total_revenue
FROM ticket_sales
WHERE year = 2020
GROUP BY movie;
```

```

    erDiagram
        Sales_fact ||--o{ Movie_dim : "FK1"
        Sales_fact ||--o{ Actor_dim : "FK2"
        Sales_fact ||--o{ Moviegoer_dim : "FK3"
        Sales_fact ||--o{ Time_dim : "FK4"
        Sales_fact ||--o{ Language_dim : "FK5"
    
```

The diagram illustrates a star schema with the following tables and attributes:

- Sales\_fact** (Fact Table):
  - PK,FK1: movie\_id
  - PK,FK2: moviegoer\_id
  - PK,FK3: time\_id
  - PK,FK4: language\_id
  - number\_of\_movie\_goers
  - revenue
- Movie\_dim** (Dimension Table):
  - PK: movie\_id
  - actor\_id
  - title
  - studio
  - director
  - release\_year
- Actor\_dim** (Dimension Table):
  - PK: actor\_id
  - actor
- Moviegoer\_dim** (Dimension Table):
  - PK: moviegoer\_id
  - name
  - gender
  - age
  - location
- Time\_dim** (Dimension Table):
  - PK: time\_id
  - day
  - month
  - year
- Language\_dim** (Dimension Table):
  - PK: language\_id
  - language

The relationships between the fact table and the dimension tables are as follows:

- Sales\_fact** to **Movie\_dim**: 1:M relationship (FK1)
- Sales\_fact** to **Actor\_dim**: 1:M relationship (FK2)
- Sales\_fact** to **Moviegoer\_dim**: 1:M relationship (FK3)
- Sales\_fact** to **Time\_dim**: 1:M relationship (FK4)
- Sales\_fact** to **Language\_dim**: 1:M relationship (FK5)

- Query to find revenue for each language of a movie  
 SELECT M.title AS Movie\_Name , L.language AS language , SUM(S.Revenue) AS Revenue  
 FROM Movie\_dim M  
 LEFT JOIN Sales\_fact S USING(movie\_id)  
 LEFT JOIN Language\_dim L USING(language\_id)  
 GROUP BY 1,2  
 ORDER BY 3 DESC;
- Query to find the most earning move for each actor  
 SELECT A.actor AS Actor, K.title AS Movie\_Name  
 FROM (  
 Select P.actor\_id , P.title , RANK() OVER(PARTITION BY P.actor\_id,P.title ORDER BY  
 P.Revenue DESC) Rank  
 FROM (  
 SELECT M.actor\_id AS actor\_id , M.title AS title , SUM(S.Revenue) AS Revenue  
 FROM Movie\_dim M  
 LEFT JOIN Sales\_fact S USING(movie\_id)  
 GROUP BY 1,2 ) P ) K  
 WHERE K.Rank = 1;

Q3)

a)

Columnar databases find application in data warehouses where enterprises funnel extensive data from various origins for Business Intelligence analysis. The efficiency of query performance in column-oriented databases is heightened due to the design that maintains data proximity, minimizing seek time.

Utilizing a columnar database addresses the constraints inherent in conventional relational databases, positioning columnar databases as the forthcoming frontier in business intelligence. In a database housing millions of records, a columnar database grants access to the most pertinent elements, enhancing query speed significantly. While traditional relational databases continue to serve as comprehensive data sources, the architecture of columnar databases simplifies the analysis of the entire dataset. The organizational structure of data in columnar databases leads to swifter outcomes and more effective analysis. Columnar databases excel at queries which involve few columns or aggregation queries against vast amounts of data.

The key observations I found from my study are:

- Numerous columnar databases employ parallel processing to horizontally expand, allowing them to manage expanding datasets and meet growing analytical requirements. This scalability proves crucial for businesses constantly accumulating data over time.
- Columnar databases can seamlessly adapt schemas by introducing new columns without disrupting existing structures or queries. This flexibility proves valuable in data warehousing, where schemas often undergo changes over time. An example of this adaptable columnar design is evident in ClickHouse, empowering organizations like Cloudflare to manage shifting data requirements.
- Materialized views empower columnar databases to pre-compute and store results for common analytical queries, significantly enhancing performance. Exasol incorporates this feature, utilized by some companies to expedite analytics.
- Columnar databases seamlessly integrate with popular business intelligence, visualization, and analytics tools, simplifying the extraction of insights from data.
- The use of columnar storage enhances the speed of query execution, especially for analytical queries, by minimizing the data processed. This efficiency stems from the storage of data in columns rather than rows, aligning seamlessly with the structure of analytical queries.
- Specialized compression techniques tailored to the storage format are frequently employed by columnar databases. This results in decreased storage needs and faster data retrieval, both vital for applications in data warehousing that manage extensive data volumes.

b)

- IBM Cognos

IBM Cognos is an analytics tool with full potential of the data with AI powered automation. The AI-powered assistant is consistently accessible – articulate our data requirements, and Cognos Analytics will craft compelling data visualizations for us. If we express a query or hypothesis for testing, AI will do the rest of the part and it will draw all the possible necessary insights.

This is a transition to the next era of Business Intelligence, featuring AI capabilities that not only deliver a precise, reliable, and comprehensive view of the business but also forecast future developments, predict outcomes, and elucidate the reasons behind them. From the IBM Cognos website I can see that they have recently launched Cognos Analytics 12.0 which is the more powerful version.

I also think that with the development of such powerful AI analytical tools the jobs in data engineer and data scientist fields will increase whereas data analyst jobs will substantially decrease. It is because most of the analytical and visualization part is performed by AI powered tools like IBM Cognos.

- Microsoft SQL Server Analysis Services (SSAS):

Analysis Services is used for corporate analytics and decision assistance as an analytical data engine (VertiPaq). It provides complex semantic data models for client applications and business reports, including Excel, Power BI, Reporting Services reports, and a range of data visualization tools.

SQL Server Analysis Services supports multidimensional models, data mining, tabular models across compatibility levels (depending on version), and Power Pivot for SharePoint. It may be installed as an on-premises server or virtual machine instance. Typically, the standard workflow involves the installation of a SQL Server Analysis Services instance, the development of a tabular or multidimensional data model, deployment of the model as a database to a server instance, processing the database to populate it with data, and subsequently granting permissions to facilitate data access. Once prepared, any client application supporting Analysis Services as a data source can access the data model.

- Oracle OLAP

A top-notch multidimensional analytic engine called Oracle OLAP is built into Oracle Database 12c. Oracle OLAP cubes produce results faster than the speed of thought by performing complex computations with straightforward SQL queries. When OLAP cubes are deployed as materialized views, the exceptional query performance can be transparently utilized to improve summary query performance against detail relational tables. Oracle OLAP enables centralized data and business rule management on a safe, scalable, and enterprise-ready platform because it is integrated into Oracle Database 12c. Analytical measures such as time-series calculations, financial models, forecasts, allocations, regressions, and more can be easily produced with Oracle OLAP. Almost

any analytical calculation need can be satisfied by simply combining hundreds of analytic functions into custom functions. A star schema design is used to represent Oracle OLAP cubes: dimension views are arranged in a constellation around the cube (or fact) view. Any reporting and analysis tool or application, including complex business intelligence solutions, SQL-based development tools, and Microsoft Excel, can easily and profitably utilize the power of Oracle OLAP thanks to this standard representation of OLAP data.

- **SAP BusinessObjects:**  
SAP BusinessObjects Business Intelligence serves as a consolidated suite designed for data reporting, visualization, and collaborative sharing. Positioned as the on-premise BI layer within SAP's Business Technology Platform, it converts data into valuable insights accessible at any time and from anywhere. This is like a BI tool with end-to-end functionalities. It is because its features include reporting, analysis, data integration, ETL, Queries, analysis tools, dashboards, Visualization, Enterprise Information Management, Security, adapting to cloud based BI solutions and also the integration with SAP and other systems. Overall it is designed to turn raw data into fruitful insights.
- **Microsoft Power BI**  
Microsoft Power BI is also a very powerful BI tool with multiple use cases. It is a robust BI tool which incorporates OLAP functionalities which helps us in enhancing data analysis and reporting capabilities. The key features of Power BI is that it allows users to connect to OLAP cubes like Microsoft SQL Server Analysis Services cubes and other multidimensional data sources. It has many functionalities like drag, drop, drill down (ex- matrix table and decomposition tree), filtering, slicing and many more. Users can also create dynamic measures and calculations using DAX in Power BI. This helps us in the creation of KPI's. Along with this there are many more features in Power BI which we can use and I have personally used them in both DATA-225 and DATA-230 subjects semester end project. I used Power BI to visualize my analysis in the form of a dashboard. According to my personal experience Power BI is easy to use and handle due to the microsoft ecosystem where if we are using Microsoft Teams then we can have sharepoint, Excel, Power BI, Microsoft Azure, Open AI all under one platform.